

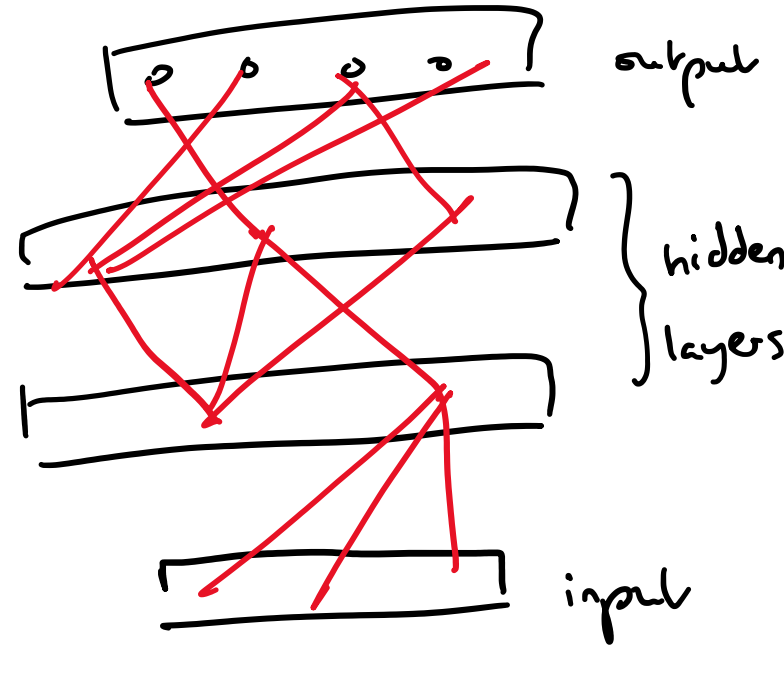
Sequence-to-sequence prediction

Both the input and the output are sequences

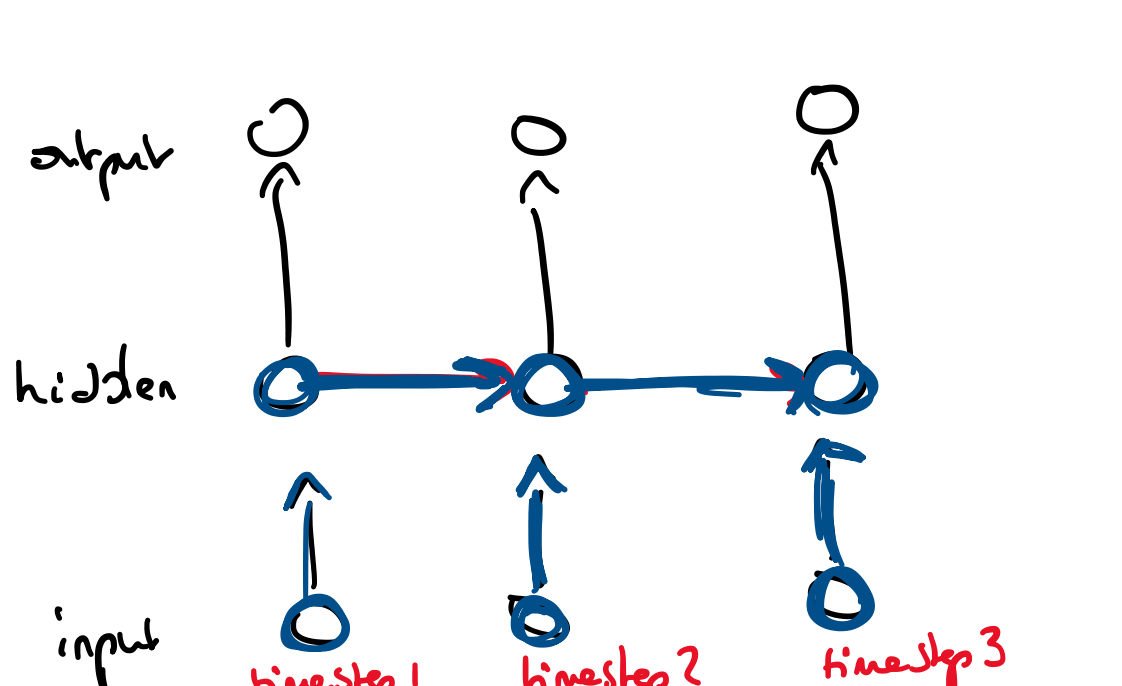
Ex applications:

- speech-to-text:
 - input: audio waveform (human speech)
 - output: text that was spoken
- alternative: text-to-speech
- caption generation:
 - input: image (not necessarily modeled as a sequence)
 - output: natural language description of the image.
- machine translation: seq-to-seq
 - input: sentence in one language (English)
 - output: _____ another language (French)

Recurrent neural networks



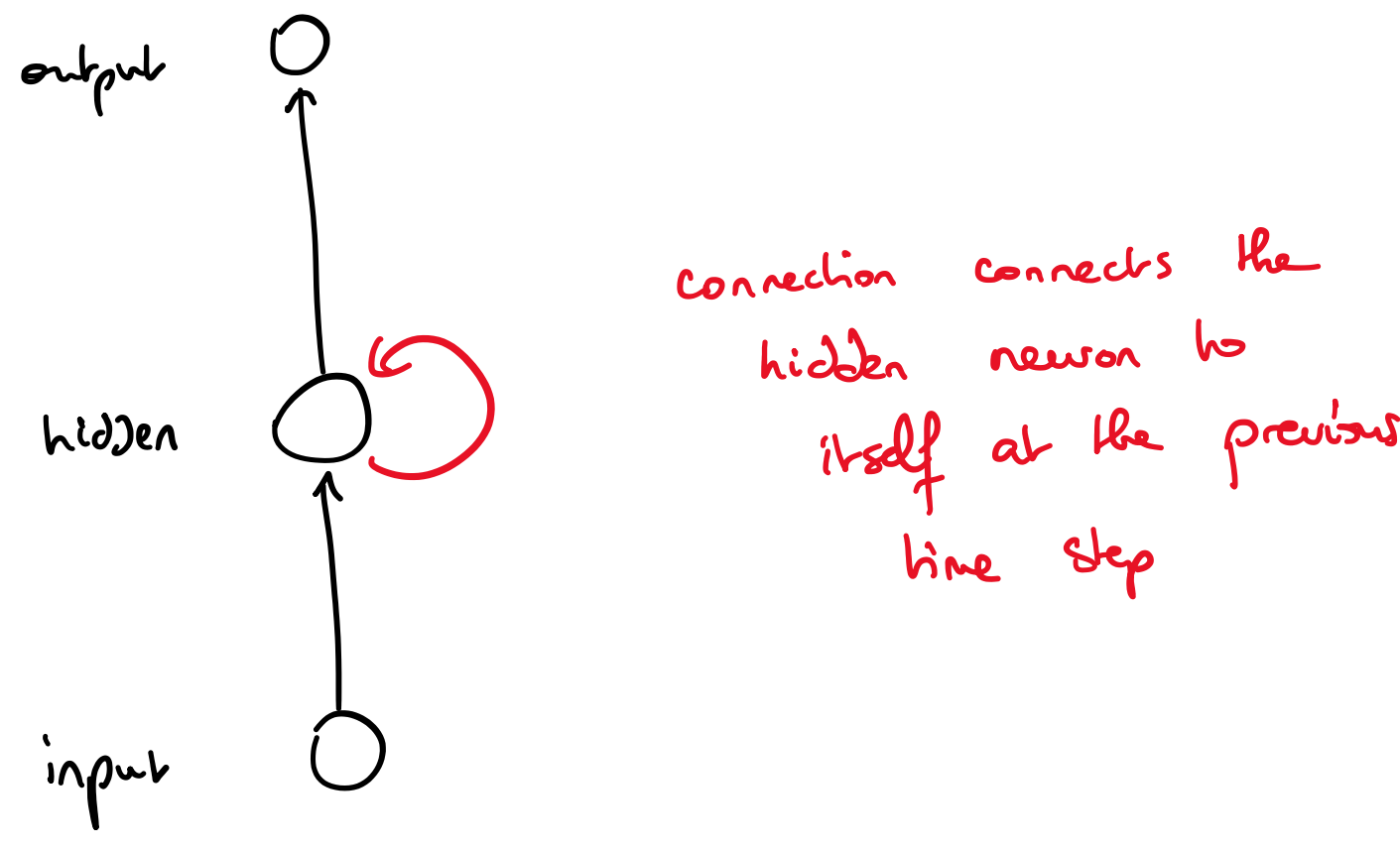
fully connected / convolutional neural networks
neurons are only connected to neurons in other layers (preceding, following)



allow connections between neurons of the same layer

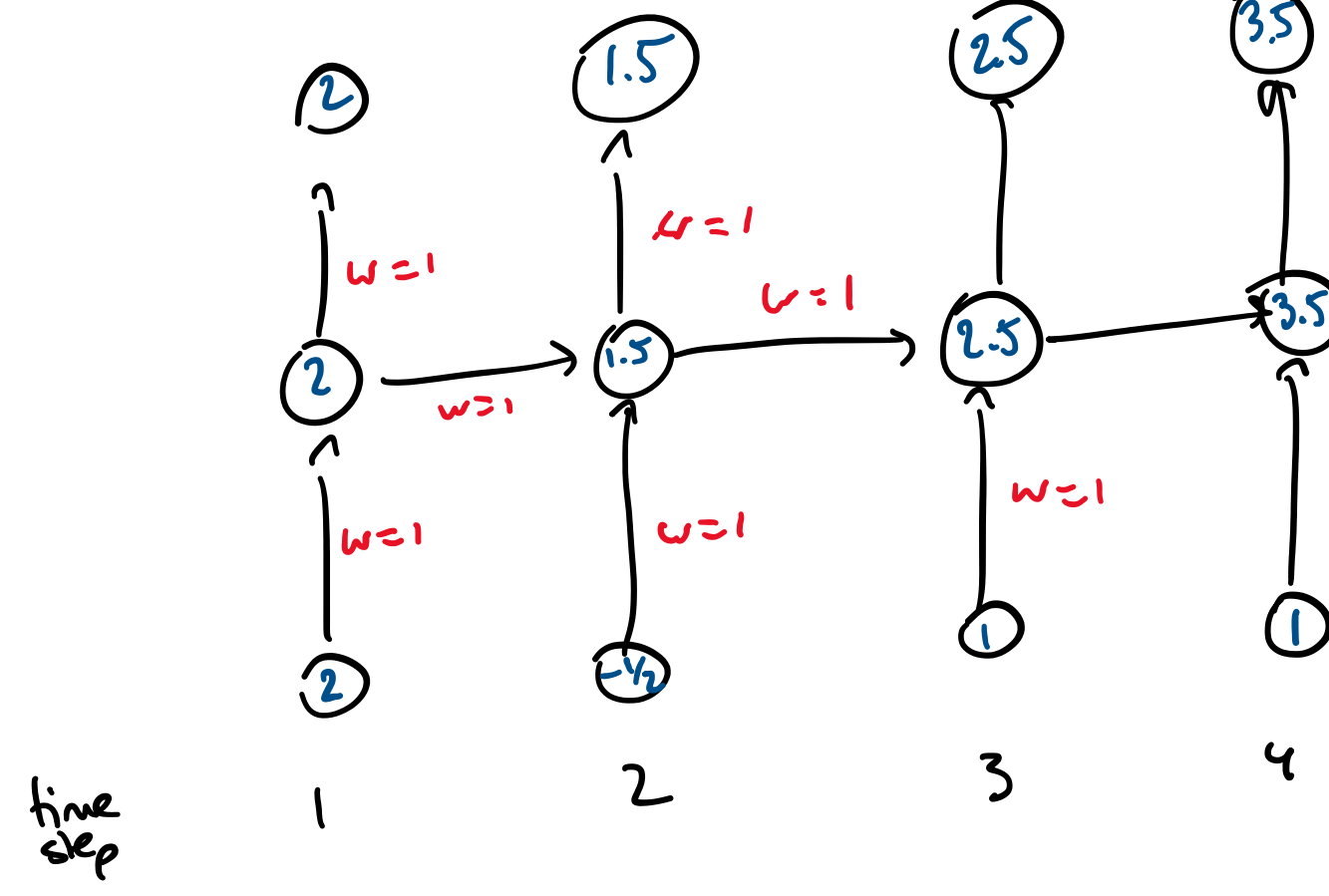
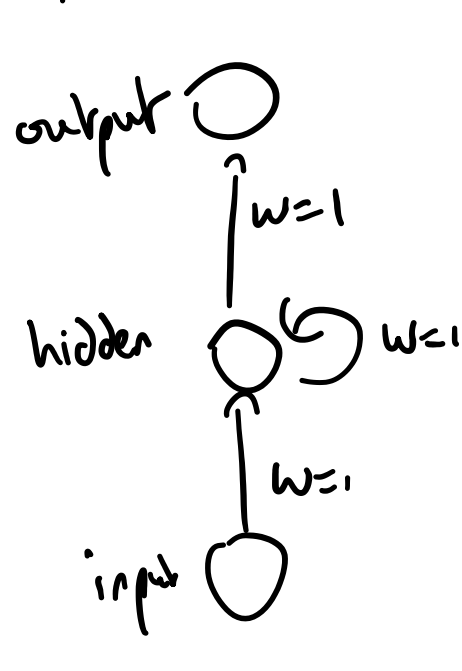
Directed graph which may contain cycles

Cycles in RNNs



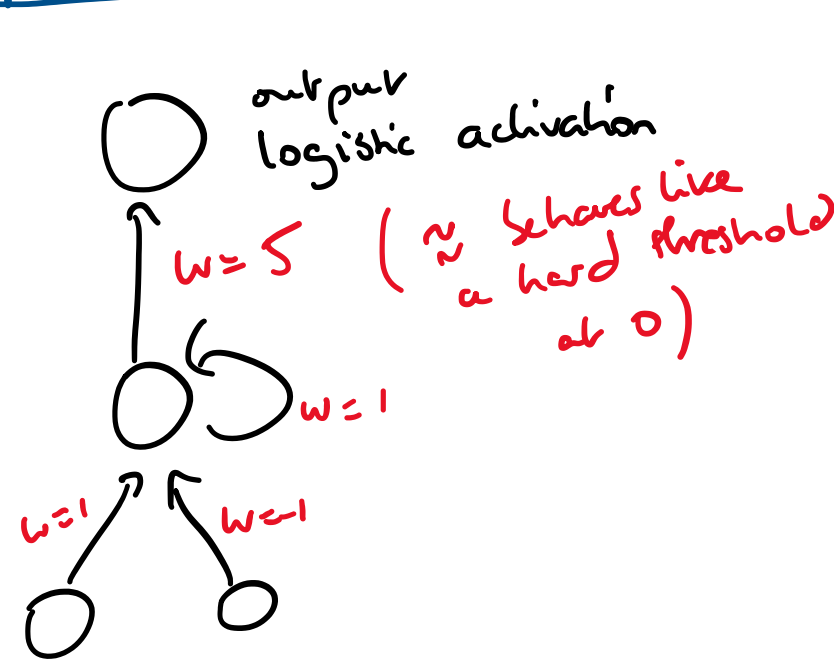
Example: RNN to sum inputs

All neurons are linear

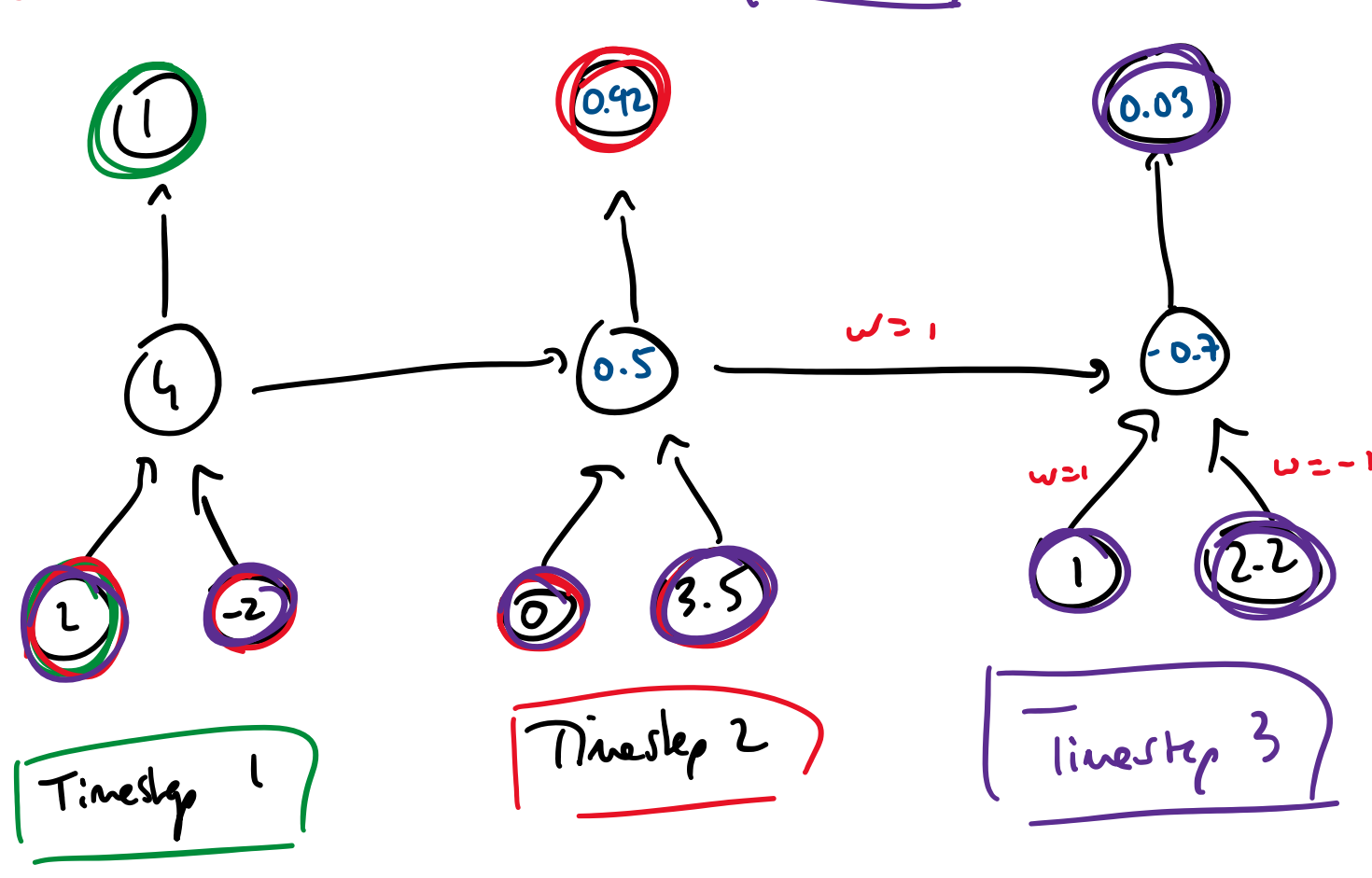


$3.5 = 1 + 1 - \frac{1}{2} + 2$

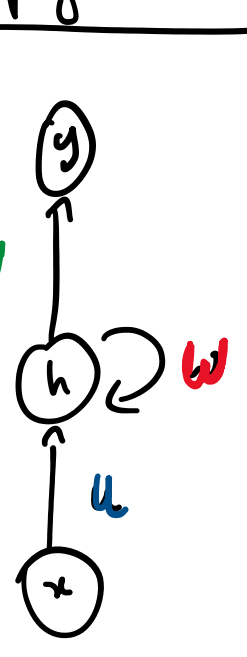
Example 2: which of the 2 inputs has a largest sum



input 1 = (2, 0, 1)
input 2 = (-2, 3.5, 2.2)

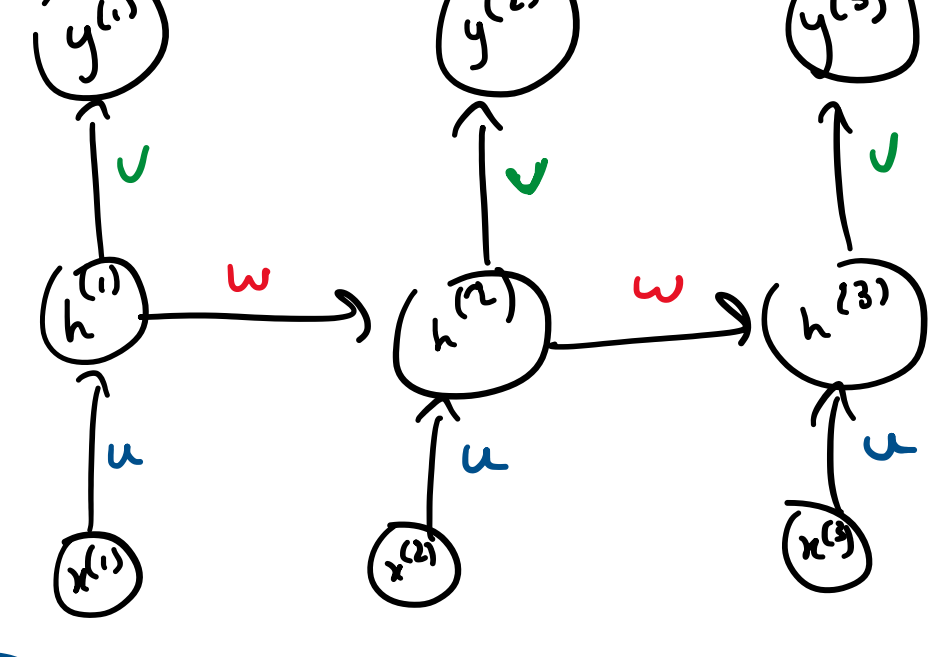


Backpropagation through time



unroll

UNROLLED COMPUTATION GRAPH



Forward pass

index of time step

$$z^{(t)} = u x^{(t)} + w h^{(t-1)}$$

activation function

$$h^{(t)} = \phi(z^{(t)})$$

result of the cycle (connection hidden-to-hidden)

$$r^{(t)} = v h^{(t)}$$

$$y^{(t)} = \phi(r^{(t)})$$

Backward pass

(we need to define the loss to derive this)

$$\bar{y}^{(t)} = \frac{\partial \mathcal{L}}{\partial y^{(t)}}$$

$$\bar{r}^{(t)} = \frac{\partial \mathcal{L}}{\partial r^{(t)}} = \frac{\partial \mathcal{L}}{\partial y^{(t)}} \frac{\partial y^{(t)}}{\partial r^{(t)}} = \bar{y}^{(t)} \phi'(r^{(t)})$$

$$\bar{h}^{(t)} = \frac{\partial \mathcal{L}}{\partial h^{(t)}} = \frac{\partial \mathcal{L}}{\partial r^{(t)}} \frac{\partial r^{(t)}}{\partial h^{(t)}} + \frac{\partial \mathcal{L}}{\partial z^{(t+1)}} \frac{\partial z^{(t+1)}}{\partial h^{(t)}} = \bar{r}^{(t)} v + \bar{z}^{(t+1)} w$$

$$\bar{z}^{(t)} = \frac{\partial \mathcal{L}}{\partial z^{(t)}} = \frac{\partial \mathcal{L}}{\partial h^{(t)}} \frac{\partial h^{(t)}}{\partial z^{(t)}} = \bar{h}^{(t)} \phi'(z^{(t)})$$

$$\bar{u} = \frac{\partial \mathcal{L}}{\partial u} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial z^{(t)}} \frac{\partial z^{(t)}}{\partial u} = \sum_{t=1}^T \bar{z}^{(t)} x^{(t)}$$

$$\bar{v} = \frac{\partial \mathcal{L}}{\partial v} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial r^{(t)}} \frac{\partial r^{(t)}}{\partial v} = \sum_{t=1}^T \bar{r}^{(t)} h^{(t)}$$

$$\bar{w} = \frac{\partial \mathcal{L}}{\partial w} = \sum_{t=2}^T \frac{\partial \mathcal{L}}{\partial z^{(t+1)}} \frac{\partial z^{(t+1)}}{\partial w} = \sum_{t=2}^T \bar{z}^{(t+1)} h^{(t)}$$

Backprop through time = forward + backward pass

Only difference w/ fully connected / convolutional neural network is summation over time steps

Train RNNs w/ algorithms like Stochastic gradient descent.

Cannot vectorize over time (because of summation)
Vectorize over training examples (minibatches) and units

⚠ RNNs often lead to vanishing/exploding gradients